# Bounded Coordinate-Descent for Biological Sequence Classification in High Dimensional Predictor Space

Georgiana Ifrim, Carsten Wiuf

{ifrim, wiuf}@birc.au.dk

*Bioinformatics Research Centre (BiRC)*

*C.F. Møllers Allé 8*

*DK-8000 Aarhus C, Denmark*

### Abstract

We present a framework for discriminative sequence classification where the learner works directly in the high dimensional predictor space of all subsequences in the training set. This is possible by employing a new coordinate-descent algorithm coupled with bounding the magnitude of the gradient for selecting discriminative subsequences fast. We characterize the loss functions for which our generic learning algorithm can be applied and present concrete implementations for logistic regression (binomial log-likelihood loss) and support vector machines (squared hinge loss). Application of our algorithm to protein remote homology detection and remote fold recognition results in performance comparable to that of state-of-the-art methods (e.g., kernel support vector machines). Unlike state-of-the-art classifiers, the resulting classification models are simply lists of weighted discriminative subsequences and can thus be interpreted and related to the biological problem.

**Keywords:** sequence classification, unrestricted-length features, wildcard matches, fast coordinate-descent, logistic regression, support vector machines, computational biology

## 1  Introduction

Many problems in biology today require accurate computational prediction of properties. For example, the primary DNA sequence is a main determinant of functional and structural protein properties, yet little is known about this relationship and we must therefore turn to computational prediction for advancing our understanding. Likewise, accurate gene and motif prediction is of crucial importance for the annotation of recently sequenced genomes. To achieve this goal we need machine learning techniques that are fast, highly scalable and preferably treat feature selection as an integral part of the learning algorithm. The latter requirement means that neither time nor expertise is invested in pre-processing the original data (e.g., for defining features) and no potentially hard to validate assumptions are made about data distribution. Recent advances in developing efficient machine learning tools such as fast logistic regression [Genkin et al., 2007, Lin et al., 2008] and support vector machines [Joachims, 2006, Hsieh et al., 2008] enable learning of classifiers in very large predictor spaces, thus reducing the need for pre-processing the data. Nevertheless, most of these techniques are typically designed to exploit the sparsity of the training set (i.e., many features occur sparsely in the training data) which holds for many applications, such as text categorization [Lin et al., 2008, Hsieh et al., 2008]. This assumption unfortunately does not hold for biological sequences, where many shorter features often occur very frequently.

We present an efficient coordinate-descent algorithm for optimizing regularized fitting of classification loss in high dimensional predictor space. Our approach does not rely on feature sparsity

assumptions. In our framework, the feature space is spanned by all subsequences present in the training set. Furthermore, the features can have a flexible number of wildcard matches, which allows us to model complex biological processes such as substitutions, insertions and deletions. The optimization proceeds coordinate-wise by iteratively selecting the feature with maximum (absolute) gradient value, following the Gauss-Southwell rule [Luenberger, 1984]. In order to select the best features fast, we provide bounds on the gradient value of any subsequence based on its prefix. This drastically reduces the search space. We discuss the tightness of the proposed bounds as well as the convergence of the algorithm.

Our learning technique is applicable to both unregularized and regularized loss functions. In this paper we show bounds for the more complex case of *elastic-net* regularized loss [Hui and Hastie, 2005]. By adding an explicit elastic-net penalty (a convex combination of $l1$ and $l2$ regularizers) to the loss function, we allow the user to directly trade-off $l1$-regularization (encouraging model sparsity) for $l2$-regularization (correcting for correlations) [Hastie et al., 2003, Hui and Hastie, 2005, Friedman et al., 2010]. As can be observed in our experiments this positively affects the prediction quality.

We present applications of our learning algorithm to protein remote homology detection and fold recognition. In order to compare to previously published results, we work with standard protein benchmarks. For homology detection we use SCOP1.59 [Jaakkola et al., 1999, Leslie et al., 2002a, Leslie and Kuang, 2004, Kuksa et al., 2008a], and for fold recognition we use the challenging dataset of Ding and Dubchak [2001]. The classification problems associated with these datasets are hard and the training data is rather small scale (2,800 and 300 sequences respectively). In order to further analyse the scalability of our technique, we present a large scale experiment on the latest version of the Silva-LSUPARC102 database (150,000 unique sequences). We compare our algorithm to state-of-the-art sequence classifiers, support vector machines with spectrum kernel [Leslie and Kuang, 2004], mismatch kernel [Leslie et al., 2002a, 2004], and the recent sparse spatial sample kernel [Kuksa et al., 2008a,b]. Besides being fast and highly accurate, our classification models are easily interpretable, an important advantage for the bioinformatics and medical communities.

The remainder of this paper proceeds as follows. In Section 2 we discuss related work. In Section 3 we present our learning algorithm. In Section 4 we discuss the experimental setup. In Section 5 we discuss the empirical results and we conclude in Section 6.

## 2    Related Work

Working in high dimensional predictor space and regularizing is statistically preferable to a two-step procedure of first reducing the dimension, then fitting a model in the reduced space [Rosset et al., 2004]. Recently, efficient regularized learning algorithms for logistic regression and support vector machines (SVM) were proposed for fitting classifiers in high dimensional predictor space [Lin et al., 2008, Hsieh et al., 2008, Perkins et al., 2003]. These algorithms are still at least linear in the number of features. For applications in which the number of features is much larger than the number of training samples, this poses a challenge for both running time and memory [Ifrim et al., 2008]. This is the case when we try to use all (unrestricted-length) subsequences in the training set as features. If the features are restricted to $k$-mers, for fixed and reasonably small $k$ (i.e., $k = 3$, or $k = 5$) and the features occur sparsely in the training instances (such as in applications related to text categorization), existing classification algorithms may still perform well [Ifrim, 2009]. However, this is not the case for applications focusing on biological sequence classification where the feature space is dense, i.e., many subsequences occur in many of the training sequences.

Recent work on efficiently computing string kernels for SVM [Leslie et al., 2002a, Leslie and Kuang, 2004, Kuksa et al., 2008a] addressed some of the the computational challenges typically

associated with this type of techniques. Nevertheless, the kernels proposed still restrict the set of features to subsequences of certain length or format (e.g., allowing a certain number of mismatches) in order to deal with the computational aspects associated with large feature spaces. In this paper we work with the unrestricted space of all subsequences in the training set and allow wildcard matches.

A popular class of models used for sequence classification is the generative classifier family, such as profile Hidden Markov Models [Eddy, 1998]. Generative models only learn from positive training examples, while discriminative techniques (such as SVM and those proposed in this paper) directly focus on separating the positive and negative training examples. Previous work has shown that discriminative approaches outperform generative approaches for sequence classification [Cheng and Baldi, 2006, Leslie et al., 2002a, Kuang et al., 2004].

In this paper we propose learning a linear classifier directly in a high dimensional feature space rather than using an implicit mapping via a string kernel as in kernel-SVM. The computational trick in our algorithm is to use coordinate descent coupled with bounding the search for the best coordinate. Our generic algorithm can be implemented for a range of classification loss functions, such as exponential loss, binomial log-likelihood loss of logistic regression, guassian loss and the squared hinge-loss of SVM [Hastie et al., 2003]. This work is an extension of the study of Ifrim et al. [2008] which focused on a simpler version of this algorithm for unregularized logistic regression applied to text categorization.

## 3 Bounded Coordinate-Descent for Sequence Classification

In this section we present our generic learning algorithm and its concrete implementation for logistic regression and support vector machines.

### 3.1 Preliminaries and Basic Notation

We first introduce the theoretical framework and some basic notation. Assume we have a training set of instance-label pairs $\{x_i, y_i\}_{i=1}^N$, with $y_i \in \{-1, +1\}$. The training instances $x_i$ are sequences, e.g., biological sequences $x_i = AGTCAACTGGAA....$, text sequences $x_i = ABCD...$, sequences of rankings $x_i = A < B < C < D < \ldots$, sequences of conjunctions of properties $x_i = A \cap B \cap C \cap D \cap \ldots$. Let $d$ be the number of distinct subsequences in the feature space. We formally represent the training sequences as binary vectors in the space of all subsequences of the training set: $x_i = (x_{i1}, \ldots x_{ij}, \ldots x_{id})^T$, $x_{ij} \in \{0, 1\}$, $i = \overline{1, N}$. Let $\beta = (\beta_1, \ldots, \beta_j, \ldots, \beta_d)$ be a parameter vector (defining a linear classifier).

The goal is to learn a mapping, also called classification model, $f : X \to \{-1, +1\}$ from the given training set such that given a new sample $x \in X$, we can predict a label $y \in \{-1, +1\}$.

Let

$$L(\beta) = \sum_{i=1}^N \xi(y_i, x_i, \beta) + CR_\alpha(\beta) \tag{1}$$

be the regularized classification loss criterion. Here, $\xi(y_i, x_i, \beta)$ is a classification loss function, $C$ is a constant and $R_\alpha(\beta)$ is a regularizer. Learning a classification model is achieved by finding the parameter vector $\beta$, that minimizes the regularized classification loss on the training set

$$\hat{\beta} = \operatorname{argmin}_\beta L(\beta).$$

In Equations (2)–(5) we list a few examples of commonly used classification loss functions.

$$Exponential\ loss : \xi(y_i, x_i, \beta) = e^{-y_i \beta^t x_i} \tag{2}$$

3

$$Binomial\ loglikelihood\ loss: \xi(y_i, x_i, \beta) = \log(1 + e^{-y_i\beta^t x_i}) \tag{3}$$

$$Squared\ hinge\ loss: \xi(y_i, x_i, \beta) = \max(1 - y_i\beta^t x_i, 0)^2 \tag{4}$$

$$Gaussian\ loss^1: \xi(y_i, x_i, \beta) = \log \frac{1}{\Phi(y_i\beta^t x_i)} \tag{5}$$

The penalty weight $C \geq 0$ controls the amount of regularization of the parameter vector $\beta$. A large $C$ means more penalty on the $\beta$ parameters. The regularization term can take various forms. The most common regularizers are the $l1$ which penalizes the sum of absolute values of $\beta_j$, and the $l2$, which penalizes the sum of squared $\beta_j$ coeficients (Equation (6)). In our work $R_\alpha(\beta)$ is the elastic-net regularizer [Hui and Hastie, 2005] (defined in Equation (7)), which is a compromise between the $l2$ ($\alpha = 0$) and the $l1$ ($\alpha = 1$) regularizers [Friedman et al., 2010].

$$l1 = \sum_{j=1}^{d} |\beta_j|, \quad l2 = \frac{1}{2}\sum_{j=1}^{d} \beta_j^2 \tag{6}$$

$$R_\alpha(\beta) = \alpha \sum_{j=1}^{d} |\beta_j| + (1-\alpha)\frac{1}{2}\sum_{j=1}^{d} \beta_j^2 \tag{7}$$

Depending on the type of regularization, increasing the weight $C$ of the penalty term results in many zero $\beta_j$ (for $l1$-regularization) or shrinking the coefficients $\beta_j$ of correlated predictors (for $l2$-regularization). The elastic-net regularizer allows balancing the two effects.

## 3.2   Generic Bounded Coordinate-Descent Algorithm

In this section we describe our learning algorithm and characterize the properties of loss functions to which it can be applied. Assume the following properties for the loss function:

1. $\xi$ depends on $y_i$, $x_i$ and $\beta$ only through the classification margin $m_i = y_i\beta^t x_i$. Note that $x_i$ is correctly classified if the margin $m_i$ is positive and the higher the margin the higher the classification confidence. We write $\xi(y, x, \beta) = \xi(m)$.

2. $\xi$ is a monotone decreasing function of the margin: $\xi'(m) \leq 0$.

3. $\xi$ is strictly convex and twice continuously differentiable.

The gradient of $L(\beta)$ (defined in Equation (1)) with respect to a coordinate $\beta_j$ is then

$$\frac{\partial L}{\partial \beta_j}(\beta) = \sum_{i=1}^{N} y_i x_{ij}\xi'(y_i\beta^t x_i) + C[\alpha\ \text{sign}(\beta_j) + (1-\alpha)\beta_j]$$

or using the margin notation

$$\frac{\partial L}{\partial \beta_j}(\beta) = \sum_{i=1}^{N} y_i x_{ij}\xi'(m_i) + CR'_\alpha(\beta_j).$$

If $R_\alpha(\beta)$ is not differentiable in $\beta = 0$ we use the left-right derivatives. We proceed minimizing $L(\beta)$ by coordinate-wise gradient descent in the feature space of all subsequences. In each iteration we

---

[1]$\Phi$ is the cdf of the standard normal distribution.

update only the coordinate corresponding to the subsequence with the largest gradient magnitude, following the Gauss-Southwell rule:

$$j = \text{argmax}_l \left| \frac{\partial L}{\partial \beta_l}(\beta) \right|$$

This results in a greedy advance towards the optimum of the objective function, without having to explicitly materialize the full space of subsequences. All we need is an efficient way to find the best coordinate (i.e., feature) in each iteration. We summarize the steps of our generic coordinate-descent algorithm in Algorithm 1.

---

**Algorithm 1** Generic coordinate-descent algorithm

---

1. Set $\beta^{(0)} = 0$

2. For t=1:T

   (a) Find $j_t = \text{argmax}_j \left| \frac{\partial L}{\partial \beta_{j_t}}(\beta^{(t-1)}) \right|$ (using Theorem 1)

   (b) Use line search to set the step length $\epsilon_t$ for coordinate $j_t$

   (c) Set $\beta_{j_t}^{(t)} = \beta_{j_t}^{(t-1)} - \epsilon_t \frac{\partial L}{\partial \beta_{j_t}}(\beta^{(t-1)})$ and $\beta_k^{(t)} = \beta_k^{(t-1)}, k \neq j_t$

---

The core part of the algorithm is a search routine (line (a) in Algorithm 1) that returns fast the best feature in each iteration. The search procedure relies on bounding the gradient value of any subsequence early on, by only looking at its prefix. This looks difficult at a first glance since the prefix may be a poor feature, while its extension can be highly discriminative. The intuition behind our bound is based on two observations: the gradient is influenced by class-wise frequency and the subsequence space has structure. More concretely, the gradient of a given feature is characterized by how many times this feature appears in the positive class versus the negative class. A feature has high gradient value if it is frequent in one class and not in the other. If we look at the gradient of a given feature class-wise, it is influenced only by the frequency of the given feature in the given class, and in turn, the frequency of the feature in that class is bounded by the frequency of its prefix. Thus, class-wise we can bound the gradient of a given subsequence based on its prefix, and then we can build a global bound from the class-wise bounds. Using these bounds we can discard large parts of the search space during the search for the best subsequence, making the whole search process efficient, both time-wise (fast running time) and space-wise (low memory).

Let $j$ be a coordinate corresponding to a given subsequence $s_j$, and $l$ be a coordinate corresponding to a super sequence of $s_j$, $s_l$, i.e. $s_j$ is a prefix of $s_l$. We write $s_j \in x_i$ to denote $x_{ij} \neq 0$.

The following theorem gives a tight upper bound on the gradient value of any subsequence.

**Theorem 1** *For any loss function $\xi$ satisfying properties 1-2 and any subsequence $s_l \supseteq s_j$, $j = 1, \ldots, d$,*

$$\left| \frac{\partial L}{\partial \beta_l}(\beta) \right| \leq \max \left\{ \left| \sum_{\{i | s_j \in x_i, y_i = +1\}} \xi'(m_i) + CR'_\alpha(\beta_l) \right|, \quad (8) \right.$$
$$\left. \left| \sum_{\{i | s_j \in x_i, y_i = -1\}} -\xi'(m_i) + CR'_\alpha(\beta_l) \right| \right\}.$$

**Proof 1** *We split the analysis in two parts, the first part focuses on deriving a bound within the negative class, and the second part focuses on the positive class. Recall that $\xi'(m) \leq 0$. Then,*

$$
\begin{aligned}
\frac{\partial L}{\partial \beta_l}(\beta) &= \sum_{\{i|s_l \in x_i\}} y_i x_{il} \xi'(m_i) + CR'_\alpha(\beta_l) \qquad (9) \\
&\leq \sum_{\{i|s_l \in x_i, y_i=-1\}} y_i x_{il} \xi'(m_i) + CR'_\alpha(\beta_l) \\
&\leq \sum_{\{i|s_j \in x_i, y_i=-1\}} y_i x_{ij} \xi'(m_i) + CR'_\alpha(\beta_l) \\
&= \sum_{\{i|s_j \in x_i, y_i=-1\}} -\xi'(m_i) + CR'_\alpha(\beta_l)
\end{aligned}
$$

*The last inequality in Equation (9) holds because of the anti-monotonicity property $\{i|s_l \in x_i, y_i = -1\} \subseteq \{i|s_j \in x_i, y_i = -1\}$. Similarly, we can show for the positive class that*

$$
\frac{\partial L}{\partial \beta_l}(\beta) \geq \sum_{\{i|s_j \in x_i, y_i=+1\}} \xi'(m_i) + CR'_\alpha(\beta_l)
$$

*Thus we obtain the lower and upper bounds*

$$
\sum_{\{i|s_j \in x_i, y_i=+1\}} \xi'(m_i) + CR'_\alpha(\beta_l) \leq \frac{\partial L}{\partial \beta_l}(\beta) \leq \sum_{\{i|s_j \in x_i, y_i=-1\}} -\xi'(m_i) + CR'_\alpha(\beta_l).
$$

*We are interested in an upper bound on the absolute magnitude of the gradient at a given coordinate, thus we write more conveniently*

$$
\left| \frac{\partial L}{\partial \beta_l}(\beta) \right| \leq \max \left\{ \left| \sum_{\{i|s_j \in x_i, y_i=+1\}} \xi'(m_i) + CR'_\alpha(\beta_l) \right|, \qquad (10) \right.
$$
$$
\left. \left| \sum_{\{i|s_j \in x_i, y_i=-1\}} -\xi'(m_i) + CR'_\alpha(\beta_l) \right| \right\}.
$$

The main property on which Theorem 1 relies on (last inequality in Equation (9)) is the anti-monotonicity: the prefix frequency is higher than that of its extension subsequence. This means that the same bound holds for *any* weight $x_j$ as long as $x_j \geq x_l$. This observation is useful for example if we want to integrate prior biological knowledge about the target problem, such as the fact that some amino acids in protein sequences are more useful than others. We could encode this information by giving weights $w_a \in (0, 1)$ to individual amino acids and still have the property that $x_j = \Pi_a w_a \geq x_l$.

As we can observe from Equation (10), without explicit regularization of the loss function, the upper bound on the gradient of a subsequence $s_l$ solely depends on the frequency of its prefix $s_j$. This means that we can decide not to expand the $s_j$ prefix further, without inspecting the longer feature $s_l$. If we use a regularizer, we have the term $\beta_l$ in the bound which depends on the longer subsequence, rather than only the prefix. Thus, we have to take care to compute the correct bound whenever $\beta_l \neq 0$. Since at start all $\beta_j$ are zero, during iterations we only need to check those features that have non-zero $\beta_l$ because they were selected in the model in previous iterations. In practice, as supported by our experiments, this part is fairly fast due to the sparsity of the final model learned by this algorithm.

The bound presented in Theorem 1 is tight since we can construct examples for which the inequality is an equality. Consider the simple case of a training set with positive examples of the type $AAAAAAA$, and negative examples of the type $BBBBBB$. Whenever the subsequence set of occurrences is the same as that of its prefix, the inequality becomes equality.

The iterates generated by Algorithm 1 converge to the optimal solution of the objective function given in Equation (1). Based on results from Luo and Tseng [1992] characterizing the convergence of iterates generated by coordinate-descent using the Gauss-Southwell rule we have:

**Theorem 2** *Let $\beta^{(t)}$ be a sequence of iterates generated by Algorithm 1. Then $\beta^{(t)}$ converges to the optimal solution of (1).*

In the next subsections we describe the specific bounds used to design learning algorithms in the style of Algorithm 1 for logistic regression and support vector machines. We choose these popular classifiers for concrete implementations because their loss functions satisfy the properties 1-3 and they are margin maximizing [Rosset et al., 2004]. The latter property typically translates into good generalization ability in practice.

### 3.3 Logistic Regression

Let

$$L(\beta) = \sum_{i=1}^{N} \log(1 + e^{-y_i \beta^t x_i}) + CR_\alpha(\beta)$$

be the elastic-net-regularized binomial log-likelihood loss [Hastie et al., 2003]. The gradient of $L(\beta)$ with respect to a coordinate $j$ at a given parameter vector $\beta$ is:

$$\frac{\partial L}{\partial \beta_j}(\beta) = \sum_{i=1}^{N} y_i x_{ij} \underbrace{\left( \frac{-1}{1 + e^{-y_i \beta^t x_i}} \right)}_{\xi'(m_i)} + CR'_\alpha(\beta_l)$$

**Corollary 1** *For binomial log-likelihood loss and any subsequence $s_l \supseteq s_j$, $j = 1, \ldots, d$,*

$$\left| \frac{\partial L}{\partial \beta_l}(\beta) \right| \leq \max \left\{ \left| \sum_{\{i | s_j \in x_i, y_i = +1\}} \frac{-1}{1 + e^{-\beta^t x_i}} + CR'_\alpha(\beta_l) \right|, \right. \tag{11}$$
$$\left. \left| \sum_{\{i | s_j \in x_i, y_i = -1\}} \frac{1}{1 + e^{\beta^t x_i}} + CR'_\alpha(\beta_l) \right| \right\}.$$

### 3.4 Support Vector Machines

Let

$$L(\beta) = \sum_{i=1}^{N} \max(1 - y_i \beta^t x_i, 0)^2 + CR_\alpha(\beta)$$

be the elastic-net-regularized squared hinge loss [Chang et al., 2008]. We can rewrite $L(\beta)$ in the equivalent form

$$L(\beta) = \sum_{\{i | 1 - y_i \beta^t x_i > 0\}} (1 - y_i \beta^t x_i)^2 + CR_\alpha(\beta)$$

7

The gradient of $L(\beta)$ with respect to a coordinate $j$ at a given parameter vector $\beta$ is:

$$\frac{\partial L}{\partial \beta_j}(\beta) = \sum_{\{i|1-y_i\beta^t x_i>0\}} y_i x_{ij} \underbrace{2(y_i\beta^t x_i - 1)}_{\xi'(m_i)} + CR'_\alpha(\beta_l)$$

**Corollary 2** *For squared hinge loss and any subsequence $s_l \supseteq s_j$, $j = 1,\ldots,d$,*

$$\left|\frac{\partial L}{\partial \beta_l}(\beta)\right| \leq \max \left\{ \left| \sum_{\{i|s_j \in x_i, 1-\beta^t x_i>0, y_i=+1\}} 2(\beta^t x_i - 1) + CR'_\alpha(\beta_l) \right|, \right. \tag{12}$$
$$\left. \left| \sum_{\{i|s_j \in x_i, 1+\beta^t x_i>0, y_i=-1\}} 2(1 + \beta^t x_i) + CR'_\alpha(\beta_l) \right| \right\}.$$

## 3.5   Algorithmic Details

In this section we give some details on using the bound of Theorem 1 to prune the search space.

Before starting the optimization iterations we build an inverted index on all the distinct unigrams in the training set, i.e., a list of document ids and positions of occurrence for each unigram. In each iteration we start the process of searching for the best feature from the level of unigrams.

We implement two strategies of prefix-expansion. The breadth-first-search (BFS) expands all unigrams to bi-grams, then all bi-grams to tri-grams, etc. For each unigram, we compute the gradient value, we keep track of the best gradient/feature seen so far and we compare the bound to the current best gradient. If the bound is lower than the current best gradient value, no subsequence starting with this unigram can improve the current optimum and we prune this part of the search space (i.e., we discard this prefix). Then we move on to the next level of expansion, and repeat this process. The depth-first-search (DFS) strategy expands the unigram to its longest subsequence until the pruning condition is met, and it then backtracks to the longest valid prefix and re-attempts to expand. Depending on the sequence tokenization used, word-level or character-level, the two strategies have different benefits. For word-level tokenization, such as in text categorization, short subsequences are more likely to be useful than long ones (e.g., phrases of 2-3 words are typically good discriminators), and therefore BFS may be a better expansion strategy. For character-level tokenization, such as for biological sequences, longer subsequences tend to be more useful than shorter ones (e.g., subsequences of length 3 or less may occur in all the sequences), so DFS is a better expansion choice. For the experiments in this paper we have used depth-first-search. In order to keep track of the occurrences of all active prefixes (e.g., subsequences that could not be pruned using Theorem 1), we expand the inverted index on-demand. The inverted index does not grow excessively due to the effectiveness of the bound. After selecting the feature with the best gradient in a given iteration we do a line search to compute the final feature weight. We use a binary search strategy for the line search: we double an initial small step size until breaking the monotonicity of the loss function (i.e., the loss function stops decreasing), after which we narrow down the search by halving the step until reaching a required precision. We stop the optimization iterations based on a convergence thereshold on the loss function.

The worst-case complexity of each iteration is $O(d + N)$, where $d$ is the size of the feature space and $N$ is the number of training samples. However, in practice our algorithm drastically prunes the search space and is thus very fast. For example, for the classification task *Scorpion-toxin like* (presented in detail in Section 5) we have observed the following behaviour. There are 23 distinct unigrams in the training set (22 amino acids and we add the wildcard as an additional unigram). Assuming we restrict the $k$-mer size to $k = 10$ (only for the scope of this example, in our experiments

$k$ is unrestricted), the total number of features is then $d = 23^{10} \sim 10^{14}$. For finding the best feature in the first iteration our algorithm checks the pruning bound 1,113 times, and prunes the space 1,104 (out of 1,113) times. The algorithm stops after 54 iterations. In the last iteration, it checks the bound 11,336 times and prunes the space 11,239 times. Thus using the bound in Theorem 1, we prune the search space to less than *a billionth of a percent* for this particular task.

For simplifying the implementation we currently prune the space using a prefix-expansion strategy. Theorem 1 however can be used for pruning all subsequences that contain a given subsequence, not necessarily in the start position as a prefix. This observation could be used to further speedup the training process.

# 4    Experiments

We implement the generic coordinate-descent algorithm in our machine learning tool SEQL (**SEQ**uence **L**earner), available from `http://www.birc.au.dk//~ifrim/seql` under the GNU GPL open source license. For now, SEQL implements the two elastic-net-regularized classifiers presented in the previous sections: logistic regression and support vector machines. In this section we describe the datasets, the techniques compared and the methodology for designing experiments.

## 4.1    Datasets

In order to compare our results to the state-of-the-art we report directly the results published in those respective papers and perform experiments on the same benchmarks.

For **protein remote homology detection** we use SCOP1.59 [Leslie et al., 2002a, Leslie and Kuang, 2004, Weston et al., 2005, Kuksa et al., 2008a]. This is an expert-curated database of protein domains organized heirarchically into folds, superfamilies and families. Protein sequences belonging to different families but the same superfamily are considered to be remote homologs in SCOP. This dataset contains 2,862 labeled sequences organized into 54 binary classification problems simulating homology detection by predicting the super-family using a hold-one-family-out strategy. No pair of sequences shares more than 95% identity. The positive training sets are quite small and the learning task is challenging. For more details on this benchmark see Jaakkola et al. [1999], Leslie and Kuang [2004], Weston et al. [2005].

For **protein remote fold recognition** we use the benchmark published by Ding and Dubchak [2001]. This dataset consists of sequences from 27 folds divided into two independent sets such that the training and test sequences share less than 35% sequence identities and within the training set, no sequences share more than 40% sequence identities. There are 311 training sequences and 383 test sequences.

In order to analyze the **scalability** of our method in a large scale experiment, we have downloaded the latest **ribosomal RNA** (rRNA) database Silva-LSUParc[2] [Pruesse et al., 2007], version 102 (released in February 2010.) The LSUParc102 database contains 180,344 rRNA sequences. After removing duplicate sequences we obtain a dataset of 150,780 unique sequences organized in 3 (one-vs-all) binary classification tasks according to the Bacteria, Archaea and Eukarya domains. The distribution by domain is dominated by Eukarya with 141,601 sequences, followed by Bacteria with 8,967 and Archaea with 212 sequences.

All datasets are available from `http://www.birc.au.dk//~ifrim/seql/data`.

---

[2]Silva-LSUParc database: http://www.arb-silva.de/documentation/background/release-102/.

## 4.2  Techniques Compared

Previous studies have shown that discriminative approaches for sequence classification (such as kernel-SVM) outperform generative approaches (such as profile HMM) by a large margin [Cheng and Baldi, 2006, Kuang et al., 2004, Leslie et al., 2002a, Liao and Noble, 2002].

We compare our algorithms to the latest techniques based on sequence kernels for SVM: the spectrum kernel, the mismatch kernel [Leslie et al., 2002a, Leslie and Kuang, 2004, Lodhi et al., 2002], and the recent sparse spatial sample kernel (SSSK) [Kuksa et al., 2008b,a]. All these methods aim at computing similarity for all pairs of sequences (the kernel matrix) in a particular feature space. Due to computational challenges, these techniques typically restrict the length and expressive power of the subsequences used as features. For example, the spectrum-$k$ kernel [Leslie et al., 2002b] implicitly compares sequences in the space of all $k$-mers, where the length $k$ of subsequence-features is a parameter of the model. The mismatch kernel [Leslie et al., 2002a] generalizes the spectrum-$k$ kernel by allowing up to $m$ mismatches or substitutions to accomodate mutations. The sparse spatial sample kernel (SSSK) further generalizes the mismatch kernel by sampling the sequences at different resolutions and comparing the resulting spectra [Kuksa et al., 2008a]. SSSK has 3 parameters, $(k, t, d)$, where $k$ is the probe size, $t$ is the number of probes and $d$ is the number of *maximum* allowed positions between the probes. Our learning algorithm uses all (unrestricted-length) subsequences in the training set as features. Furthermore, we also allow mismatches or so called wildcard matches, by a parameter that controls the *maximum* number of consecutive wildcards allowed. This allows us to model complex biological processes such as substitutions, insertions and deletions. Figure 1 gives examples for the types of features implicitly used by the above described kernels, as compared to the features used by our technique.

One advantage of kernel techniques is their ability to integrate unlabeled data during kernel computation to relax labeled data requirements. However, computing all-pair similarities for large datasets (resulting from the addition of unlabeled sequences) remains computationally very challenging for kernel methods both time and memory-wise. For example, we couldn't apply any of the above mentioned kernel methods on our large dataset (150,000 sequences) since this would require more than 90GByte memory.

For completion we also report the results of a technique published by Liao and Noble [2002] which applies SVM using an empirical kernel map based on pairwise Smith-Waterman sequence alignment scores (SVM-pairwise in Table1).

We use SEQL-LR and SEQL-SVM to denote the logistic regression and support vector machines implementations of our generic learning algorithm.

## 4.3  Methodology

Previous studies on biological sequence classification have extensively used the AUC (area under the ROC curve) and AUC50 for evaluation [Fawcett, 2004, Leslie and Kuang, 2004, Sonego et al., 2008]. AUC describes the ranking of prediction scores rather then being dependent on a fixed classification threshold. The ROC curve is obtained by plotting the fraction of true positives versus the fraction of false positives for a binary classifier as its discrimination threshold is varied [Fawcett, 2004]. A common aggregate measure is to report the area under the ROC curve (AUC) where an area of 1 represents a perfect ranking of all positives above all negatives and an area close to 0.5 represents a random classifier. The AUC50 focuses on top ranked examples and is defined as the normalized area under the ROC curve computed for up to 50 true negatives [Gribskov and Robinson, 1996]. It is typically used to evaluate classifiers on datasets where the number of positives is much lower than the number of negatives. While for a complete test set the AUC is betwen 0.5 and 1, the AUC values for truncated top lists are between 0 and 1. The shorter the top list is, the smaller the AUC
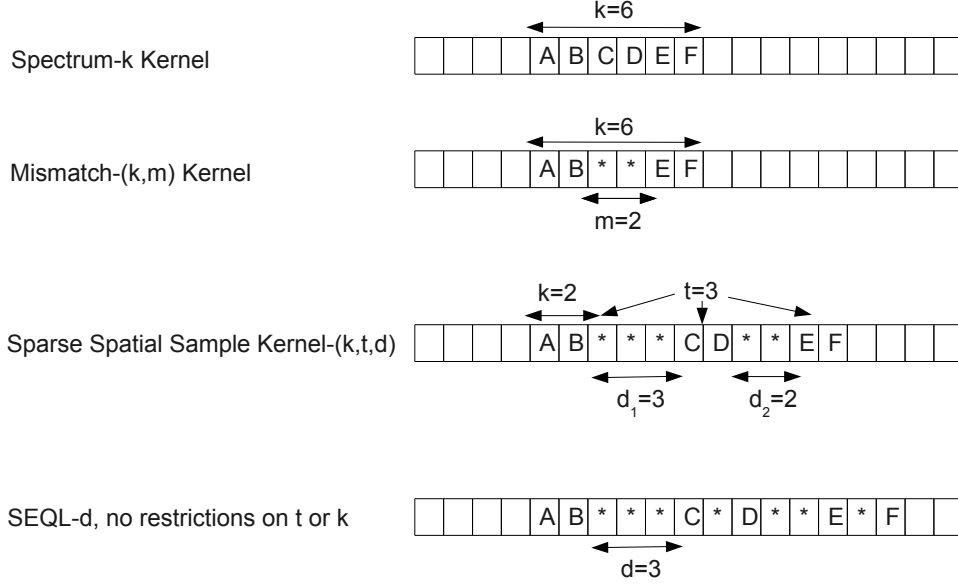
Figure 1: Features employed by state-of-the-art kernels versus our method SEQL. **Spectrum-k Kernel**, **k** is the length of features. **Mismatch-(k,m) Kernel**, **k** is the length of features, **m** is the maximum number of mismatches. **Sparse Spatial Sample-(k,t,d) Kernel**, **k** is the probe size, **t** is the number of probes and **d** is the number of maximum allowed positions between the probes. SEQL-**d**, **d** is the maximum number of consecutive wildcards.

values [Sonego et al., 2008]. In some studies instead of AUC or AUC50 only the balanced-error-rate is reported. To compare our results to state-of-the-art, we also show the balanced-error (Equation (13)) which measures the average of the errors on each class for a fixed classification threshold. It can equivalently be interpreted in terms of Specificity and Sensitivity [Levner et al., 2006].

$$BER = \frac{1}{2} \left( \underbrace{\frac{FN}{TP+FN}}_{1-\text{Sensitivity}} + \underbrace{\frac{FP}{FP+TN}}_{1-\text{Specificity}} \right) \qquad (13)$$

The benchmarks for remote homology detection and fold recognition have pre-defined training and test splits. We use the same data for experiments for all methods compared, and for the prior techniques *we report directly the results published in those respective papers*. For the large scale experiment on ribosomal RNA, we show 5-fold cross-validation results on the full dataset. Since the prior techniques required more than the available memory, for the large dataset we only report results using our method.

All experiments were run on a Linux machine with 2.5GByte memory and 2.8GHz Intel CPU.

# 5   Results and Discussion

In this section we present the results of applying SEQL-LR and SEQL-SVM to protein remote homology detection and fold recognition. Furthermore, to evaluate the scalability of our approach, we present

a large-scale experiment on the latest release of the Silva-LSUParc database [Pruesse et al., 2007].

## 5.1 Comparison to State-of-the-Art

We first compare the performance of SEQL-LR and SEQL-SVM to that of kernel-SVM techniques on two protein classification tasks.

### 5.1.1 Protein remote homology detection

In Table 1 we show results on the SCOP1.59 benchmark for all compared methods. The numbers in brackets next to the name of each method refer to the explicit parameters of these methods as discussed in the previous section, e.g., length of subsequences used as features or maximum number of wildcards. The results are averaged over all 54 binary classification tasks which correspond to superfamilies.

| Method | AUC | AUC50 | BER |
|---|---|---|---|
| SPECTRUM-(2) | 0.8581 | 0.3583 | - |
| SPECTRUM-(3) | 0.8723 | 0.4037 | - |
| MISMATCH-(5,1) | 0.8749 | 0.4167 | - |
| SSSK(1,2,5) | 0.8901 | 0.4629 | - |
| SSSK(1,3,3) | 0.9148 | 0.5118 | - |
| SVM-pairwise | 0.8930 | 0.4340 | - |
| SEQL-LR(5) | **0.9106** | **0.5185** | **0.4418** |
| SEQL-SVM(5) | **0.9214** | **0.5213** | **0.4383** |

Table 1: **Remote homology detection on the Scop1.59 dataset.** The average AUC, AUC50 and BER scores over the 54 target superfamilies. Results for kernel-SVM methods cited directly from Leslie and Kuang [2004], Kuksa et al. [2008a].

For SEQL-LR and SEQL-SVM we set the maximum allowed number of wildcards to 5. We varied the regularization parameters for the elastic net penalty. We show the most accurate results over the set of parameter values tried. We observe that even if in this benchmark the positive training sets are quite small (minimum number of samples is 2 and maximum is 95) our techniques outperform state-of-the-art methods. We believe this is due to the additional flexibility of features used for learning classifiers. Our approach does not restrict the length of features, but rather lets the data drive the classifier's decisions as to which subsequences are most discriminative. We also note that SEQL-SVM has slightly better AUC and AUC50 than SEQL-LR . Learning SEQL classifiers on this dataset took a few seconds per topic and about 80MByte memory. The computational requirements of SEQL are mainly influenced by the maximum allowed number of wildcards (set to 5 in this experiment).

Figure 2 shows a graphic comparison of some of the methods (for which we had access to published per class AUC50 scores) by plotting the total number of superfamilies above an AUC50 score threshold as a function of the threshold value. The threshold is changed so as to generate the next AUC50 value in the ranked list of AUC50 scores of the respective method. The numerical integral of this cumulative AUC50 curve is the arithmetic average of the AUC50 values shown in Table 1 [Sonego et al., 2008]. SVM with sparse-sample-spatial kernel (SSSK-SVM) behaves best among the prior techniques. We observe that SEQL-LR and SEQL-SVM are comparable to SSSK-SVM with some gain on families where SSSK-SVM has either low or good quality (AUC50 zero or above 0.4). Even though the features employed by SSSK-SVM are quite flexible (see Figure 1), we still gain some performance by imposing even less restrictions on features.

Figure 2: **Comparison of mismatch-SVM, SSSK-SVM, SEQL-LR and SEQL-SVM on SCOP1.59.** The graph plots the total number of superfamilies for which a given method exceeds an AUC50 score threshold.

Kuksa et al. [2008a] analyzed the biological relevance of features learned by sssk-svm taking the *Scorpion toxin-like* superfamily as an example. This classification task has 16 members of this superfamily as positive training examples, 1067 non-members as negative examples and the *short-chain scorpion toxin* family as a test case. sssk-svm(1,2,5) achieved an AUC50 of 0.7661 for this task [Kuksa et al., 2008a]. In their work Kuksa et al. [2008a] present the schematic representation of the *short-chain scorpion toxin* family obtained from PROSITE [Hulo et al., 2006]. The PROSITE database consists of a large collection of manually-curated biologically meaningful signatures that are described as patterns or profiles [Hulo et al., 2006]. In Figure 3 we show the PROSITE representation. The scheme shows that these type of toxins contain six conserved cysteines (C) residues



Figure 3: Schematic representation of the *short-chain scorpion toxins* family from the PROSITE database.

involved in disulfide bonds. PROSITE also lists a consensus pattern present in all the members of this family shown as well in Figure 3 (marked with 'X'). We focus here on the same superfamily and look at the features learned by seql-svm(5). In Table 2 we show the top-10 positive and negative features selected by seql-svm on this superfamily.

seql-svm(5) has an AUC50 of 0.8847 on this task. Kuksa et al. [2008a] analyzed the support vectors of their method and observed that the feature C***C had the highest weight. Similarly, seql-svm(5) selects the feature C***C in top-10, but ranks features such as N*C***C and C*C higher, which boosts its AUC50 score. In order to understand the effect of these features on the

| Positive Features | | Negative Features | |
| --- | --- | --- | --- |
| Weight | Feature | Weight | Feature |
| 0.1551 | SG*C | -0.1449 | G |
| 0.1518 | GYC | -0.0783 | CP |
| 0.1388 | N**C***C | -0.0754 | T |
| 0.0885 | C*C | -0.0541 | R****L |
| 0.0816 | C*****A****C | -0.0502 | S*R |
| 0.0657 | G***G*C | -0.0462 | C***T |
| 0.0654 | C***C***G | -0.0428 | A**L |
| 0.0484 | C***C | -0.0405 | S*****K |
| 0.0479 | C*****C***C | -0.0400 | C**G |
| 0.0394 | K**C | -0.0382 | E |

Table 2: Top-10 positive and negative SEQL-SVM features on the SCOP.1.59 Scorpion-toxin-like superfamily.

AUC50 score, in Table 3 we take a closer look at the training/test statistics of some of them.

| | Training (superfamily) | | Test (family) | |
| --- | --- | --- | --- | --- |
| Feature | Positive | Negative | Positive | Negative |
| 1. SG*C | 62.5% (10/16) | 1.5% (17/1067) | 13% (3/23) | 1.2% (21/1661) |
| 2. GYC | 62.5% (10/16) | 0.5% (6/1067) | 0.4% (1/23) | 0.09% (15/1661) |
| 3. N**C***C | 62.5% (10/16) | 0.02% (3/1067) | 0.4% (1/23) | 0.03% (6/1661) |
| 4. C*C | 100% (16/16) | 5.7% (61/1067) | 100% (23/23) | 4.4% (74/1661) |
| 8. C***C | 100% (16/16) | 6.3% (68/1067) | 100% (23/23) | 6.1% (102/1661) |

Table 3: Quality of the top ranked positive SEQL-SVM(5) features for the Scorpion-toxin-like superfamily.
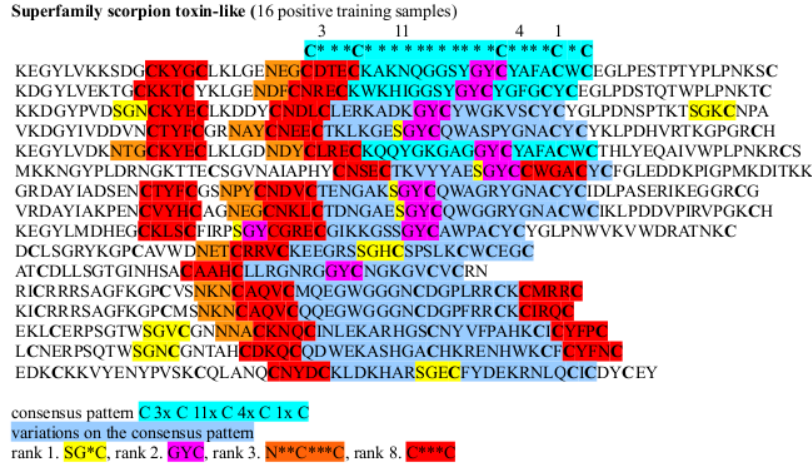


Figure 4: Position of top SEQL-SVM features relative to the PROSITE consensus pattern.

We show the percentage of occurrence followed by absolute numbers in the brackets. The notation 10/16 means that the respective feature occurs in 10 samples out of 16. Since the dataset is very unbalanced, percentages are a better indicator of feature quality. We observe that although of high

quality, the feature C***C also occurs in quite a few negatives and it is less useful than C*C for example. Note that the final ranking does not necessarily correspond to the order of selection during optimization iterations. For example the feature C*C was selected before SG*C in the optimization iterations, but due to the specific loss function, the repeated selection of SG*C lead to a slightly higher weight for this feature in the final model.

In order to asses the importance of these automatically selected features relative to the manually-curated pattern, we show in Figure 4 their exact location in the training sequences. We note that most of these features are part of the consensus pattern.

### 5.1.2 Protein remote fold recognition

In this section we show an application of SEQL-LR and SEQL-SVM to protein remote fold recognition on the dataset published by Ding and Dubchak [2001]. They proposed a technique based on building an SVM classifier using features deemed biologically relevant for this problem, such as percentage composition of the 20 amino acids, predicted secondary structure, polarity. In Table 4 we show results comparing the technique of Ding and Dubchak [2001], denoted by SVM(D&D), with published results of kernel-SVM methods and with our SEQL-LR and SEQL-SVM algorithms. Recall that lower BER means better classification quality.

| Method | AUC | AUC50 | BER |
|---|---|---|---|
| SVM (D&D) | - | - | 0.5650 |
| MISMATCH(5,1) | - | - | 0.5322 |
| SSSK(1,2,5) | - | - | 0.4619 |
| SSSK(1,3,3) | - | - | 0.4499 |
| SEQL-LR(5) | **0.7705** | **0.3519** | **0.4407** |
| SEQL-SVM(5) | **0.7848** | **0.3317** | **0.4289** |

Table 4: **Remote fold recognition on the D&D dataset.** The average AUC, AUC50 and BER scores over the 27 target folds. Results for SVM(D&D) and for the other kernel-SVM methods cited directly from Ding and Dubchak [2001], Kuksa et al. [2008a].

Even if we don't use any domain-specific knowledge in this experiment, we observe that our techniques have better classification quality than the method of Ding and Dubchak [2001] which relies on deep biological insight. We believe this advantage results from treating feature selection as a part of the learning algorithm and allowing for a large degree of flexibility in the features. SEQL-LR and SEQL-SVM outperform the state-of-the-art techniques in terms of balanced-error-rate on this benchmark.

## 5.2 Large-Scale Experiment

In this section we analyze the scalability of our method on a ribosomal RNA (rRNA) domain-prediction task.

## 5.3 Ribosomal RNA Domain-Prediction

The dataset used for this task contains 150,780 unique rRNA sequences. It is estimated that based on the new capacity for cheap and rapid sequencing there is a steady flow of about 10,000 rRNA sequences per month into the public sequence databases [Pruesse et al., 2007]. Furthermore, many sequences are derived from cultivation independent biodiversity surveys, which rely on rapid pattern- or clone-based approaches that often generate partial rRNA sequences [Pruesse et al., 2007].

We show the average results of SEQL-LR and SEQL-SVM over 5-fold cross-validation splits for fixed parameter values. We set the amount of regularization $C = 1.0$, we balance $l1$ and $l2$ regularization by setting $\alpha = 0.5$ and allow no wildcards by setting the maximum number of consecutive wildcards $d = 0$.

In Table 5 we show classification quality, running time and memory resources required by SEQL-LR and SEQL-SVM . We note that both SEQL-based techniques have high classification quality.

| Method | AUC | AUC50 | BER | Time | Memory |
|---|---|---|---|---|---|
| SEQL-LR | 0.9981 | 0.9722 | 0.0088 | 30 min | 2GByte |
| SEQL-SVM | 0.9987 | 0.9653 | 0.0083 | 20 min | 2GByte |

Table 5: Ribosomal RNA domain-prediction on Silva-LSUParc102.

Additionally, they only take about half-an-hour running-time and a reasonable 2GByte memory instead of 90GByte required by a kernel matrix computation. Instead of building sequence classifiers by costly multiple sequence alignment as currently done for sequence retrieval in this database, we could directly learn domain SEQL-classifiers from the original input sequences in the respective domains.

# 6    Conclusion

In this paper we have presented a new learning algorithm for sequence classification in high dimensional predictor space. The algorithm has at its core a coordinate-wise gradient descent strategy coupled with bounding the magnitude of the gradient to retrieve high quality features fast. We have characterized the loss functions to which our generic learning algorithm can be applied and have presented concrete implementations for logistic regression and support vector machines. Based on applications to protein remote homology detection and fold recognition, as well as large-scale domain-prediction for ribosomal RNA, we have observed our techniques are comparable to state-of-the-art in terms of classification quality. In addition, the techniques presented are highly scalable and the resulting classification models can easily be interpreted and connected to biologically relevant facts.

# Appendix A. Parameter Options for SEQL

In this appendix we give a list of SEQL parameters the user can set in order to influence the outcome classification model. SEQL is available from `http://www.birc.au.dk/~ifrim/seql`.

⟨**-o objective**⟩ Objective function. Choice between logistic regression and support vector machines. By default set to logistic regression.

⟨**-g maxgap**⟩ Maximum number of consecutive gaps or wildcards allowed in a feature, e.g., A**B, is a feature of size 4 with any 2 characters from the input alphabet in the middle. By default set to 0.

⟨**-C regularizer_value**⟩ Value of the regularization parameter. By default set to 1.

⟨**-a alpha**⟩ Weight of l1 vs l2 regularizer for the elastic-net penalty. By default set to 0.5.

⟨**-l minpat**⟩ Threshold on the minimum length of any feature. By default set to 1.

⟨**-L maxpat**⟩ Threshold on the maximum length of any feature. By default the maximum length is unrestricted, i.e., at most as long as the longest sequence in the training set.

⟨**-m minsup**⟩ Threshold on the minimum support of features, i.e., number of sequences containing a given feature. By default set to 1.

⟨**-n token_type**⟩ Word or character-level token. Words are delimited by white spaces. By default set to character-level tokens.

⟨**-r traversal_strategy**⟩ Breadth First Search or Depth First Search traversal of the search tree. By default set to DFS.

⟨**-c convergence_threshold**⟩ Stopping threshold based on change in aggregated score predictions. By default set to 0.005.

⟨**-T maxitr**⟩ Number of optimization iterations. By default set to the maximum between 5,000 and the number of iterations resulting by using a convergence threshold on the aggregated change in score predictions.

⟨**-v verbosity**⟩ Amount of printed detail about the training of the classifier. By default set to 1 (light profiling information).

# References

K.-W. Chang, C.-J. Hsieh, and C.-J. Lin. Coordinate descent method for large-scale l2-loss linear support vector machines. *Journal of Machine Learning Research*, 9:1369–1398, 2008.

J. Cheng and P. Baldi. A machine learning information retrieval approach to protein fold recognition. *Bioinformatics*, 22(12):1456–1463, 2006.

C. H.Q. Ding and I. Dubchak. Multi-class protein fold recognition using support vector machines and neural networks . *Bioinformatics*, 17(4):349–358, 2001.

S. R. Eddy. Profile hidden markov models. *Bioinformatics*, 14(9):755–763, October 1998.

T. Fawcett. Roc graphs: Notes and practical considerations for researchers, 2004.

Jerome H. Friedman, Trevor Hastie, and Rob Tibshirani. Regularization paths for generalized linear models via coordinate descent. *Journal of Statistical Software*, 33(1):1–22, 2 2010.

A. Genkin, D. Lewis, , and D. Madigan. Large-scale bayesian logistic regression for text categorization. *Technometrics*, 49(3):291–304, 2007.

M. Gribskov and N. L. Robinson. Use of receiver operating characteristic (roc) analysis to evaluate sequence matching. *Computers and Chemistry*, 20(1):25 – 33, 1996.

T. Hastie, R. Tibshirani, and J. Friedman. *The Elements of Statistical Learning*. Springer Series in Statistics, 2003.

C.-J. Hsieh, K.-W. Chang, C.-J. Lin, S. S. Keerthi, and Sellamanickam Sundararajan. A dual coordinate descent method for large-scale linear SVM. In *Proceedings of the International Conference on Machine Learning*, 2008.

Z. Hui and T Hastie. Regularization and variable selection via the elastic net. *Journal Of The Royal Statistical Society Series B*, 67(2):301–320, 2005.

N. Hulo, A. Bairoch, V. Bulliard, L. Cerutti, E. De Castro, P. S. Langendijk-Genevaux, M. Pagni, and C. J. A. Sigrist. The prosite database. *Nucleic Acids Research*, 34(1):227–230, 2006.

G. Ifrim. *Statistical Learning Techniques for Text Categorization with Sparse Labeled Data*. PhD thesis, Universität des Saarlandes, February 2009.

G. Ifrim, G. Bakir, and G. Weikum. Fast logistic regression for text categorization with variable-length n-grams. In *Proceedings of the ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, pages 354–362, New York, NY, USA, 2008. ACM.

T. Jaakkola, M. Diekhans, and D. Haussler. Using the fisher kernel method to detect remote protein homologies. In *Proceedings of the International Conference on Intelligent Systems for Molecular Biology*, pages 149–158. AAAI Press, 1999.

T. Joachims. Training linear SVMs in linear time. In *Proceedings of the ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, pages 217–226, New York, NY, 2006. ACM Press.

R. Kuang, E. Ie, K. Wang, K. Wang, M. Siddiqi, Y. Freund, and C. Leslie. Profile-based string kernels for remote homology detection and motif extraction. In *Proceedings of IEEE Computational Systems Bioinformatics Conference*, pages 152–160, Washington, DC, USA, 2004. IEEE Computer Society.

P. Kuksa, P. H. Huang, and V. Pavlovic. A fast, semi-supervised learning method for protein sequence classification. In *International Workshop on Data Mining in Bioinformatics*, pages 29–37, 2008a.

P. Kuksa, P. H. Huang, and V. Pavlovic. Scalable algorithms for string kernels with inexact matching. In D. Koller, D. Schuurmans, Y. Bengio, and L. Bottou, editors, *Neural Information Processing Systems*, pages 881–888. 2008b.

C. Leslie and R. Kuang. Fast string kernels using inexact matching for protein sequences. *Journal of Machine Learning Research*, 5:1435–1455, 2004.

C. Leslie, E. Eskin, and W. Noble. Mismatch string kernels for svm protein classification. In *Neural Information Processing Systems*, pages 1441–1448, 2002a.

C. Leslie, E. Eskin, and W.S. Noble. The spectrum kernel: a string kernel for svm protein classification. In *Proceedings of the Pacific Biocomputing Symposium*, pages 564–575, 2002b.

C. Leslie, E. Eskin, A. Cohen, J. Weston, and W. S. Noble. Mismatch string kernels for discriminative protein classification. *Bioinformatics*, 20(4):467–476, 2004.

I. Levner, V. Bulitko, and G. Lin. Feature extraction for classification of proteomic mass spectra: A comparative study. *Studies in Fuzziness and Soft Computing*, 207:607–624, 2006.

L. Liao and W. S. Noble. Combining pairwise sequence similarity and support vector machines for remote protein homology detection. In *RECOMB '02: Proceedings of the International Conference on Computational Biology*, pages 225–232, New York, NY, USA, 2002.

C.-J. Lin, R. C. Weng, and S. S. Keerthi. Trust region Newton method for large-scale logistic regression. *Journal of Machine Learning Research*, 9:627–650, 2008.

H. Lodhi, C. Saunders, J. Shawe-Taylor, N. Cristianini, and C. Watkins. Text classification using string kernels. *Journal of Machine Learning Research*, 2:419–444, 2002.

D. Luenberger. *Linear and Nonlinear Programming.* Addison-Wesley Publishing Co.: Reading Mass, 1984.

Z. Q. Luo and P. Tseng. On the convergence of the coordinate descent method for convex differentiable minimization. *Journal of Optimization Theory and Applications*, 72(1):7–35, 1992.

S. Perkins, K. Lacker, and J. Theiler. Grafting: fast, incremental feature selection by gradient descent in function space. *Journal of Machine Learning Research*, 3:1333–1356, 2003.

E. Pruesse, C. Quast, K. Knittel, B. M. Fuchs, W. Ludwig, J. Peplies, and F. O. Glockner. SILVA: a comprehensive online resource for quality checked and aligned ribosomal RNA sequence data compatible with ARB. *Nucleic Acids Research*, 35(21):7188–7196, 2007.

S. Rosset, J. Zhu, and T. Hastie. Boosting as a regularized path to a maximum margin classifier. *Journal of Machine Learning Research*, 5:941–973, 2004.

P. Sonego, A. Kocsor, and S. Pongor. ROC analysis: applications to the classification of biological sequences and 3D structures. *Briefings in Bioinformatics*, 9(3):198–209, 2008.

J. Weston, C. Leslie, E. Ie, D. Zhou, A. Elisseeff, and W. S. Noble. Semi-supervised protein classification using cluster kernels. *Bioinformatics*, 21(15):3241–3247, 2005.